



PHENIX, a high throughput structure determination tool for single-wavelength anomalous diffraction data

S Selvanayagam,¹ D Velmurugan,^{1*} T Yamane² and A Suzuki²

¹Department of Crystallography and Biophysics, University of Madras,
Guindy Campus, Chennai - 600 025, India

²Department of Biotechnology and Biomaterial Science, Graduate School of Engineering,
Nagoya University, Furo-cho, Chikusa-ku, Nagoya 464-8603, Japan

E-mail: dvelu@yahoo.com

Received 28 April 2006, accepted 28 August 2006

In Structural Genomics, one is interested in determining the structure in the fastest way to understand new folds and this has opened 'High Throughput Crystallography'. This High Throughput structure determination requires automation to reduce the obstacles related to macromolecular crystallography. PHENIX is a new software package for the automation of crystallographic structure solution of macromolecules. The algorithm allows one to proceed from reduced intensity data to a refined molecular model and to facilitate structure solution for both the expert and novice crystallographers. In cases where PHENIX builds an incomplete model of less residues with or without side chains, ARP/wARP can be made use of further with this as input. Attempts are here made in extending the use of PHENIX for the high throughput structure elucidation of proteins of approximately 44 kDa molecular weight with lab source CuK α and CrK α anomalous scattering data sets corresponding to 1.7 and 2.1 Å resolution respectively. One manganese position initially located by PHENIX is enough to build the entire structure using the CuK α data whereas 12 heavy atoms (11 heavy atoms) are required in case of CrK α data to build the entire structure.

ds: SAD, PHENIX, glucose

CS: N50, 61, 10, N4, 87, 14, E6

Introduction

Macromolecular crystallography has undergone tremendous progress in the last decade. The database of known sequences is growing at an exponential rate and has already resulted in complete genomic sequences for several organisms. The functional significance of coding regions of the genomes can be ascertained when the sequence information can be interpreted in the light of known structures. But only a fraction of the structural information is available at present. Due to the increasing rate of protein structure determination, the structural efforts of the sequencing projects have been accelerated. But even then there remains a huge sequence data gap in which three dimensional structural information is unknown.

The technologies needed for the steps of structure determination are becoming reliable and powerful enough that

they can be linked together into an automatic sequence that can yield a structure in an automatic or nearly automatic fashion. The high-throughput structure determination will require automation to reduce the obstacles related to human intervention. Several projects are underway worldwide to automate the process of sample preparation, crystallization and data acquisition. The fully automated structure determination has supplied much of the impetus for the vision of large-scale structure determination in structural genomics. In future, macromolecular crystallography would have matured considerably to such an extent that most protein crystal structures will be solved in an almost completely automatic way. This is mostly important for researchers in pharmaceutical companies.

With recent progress in data collection techniques and the current trend towards high-throughput structure determination, the single-wavelength anomalous diffraction (SAD) approach

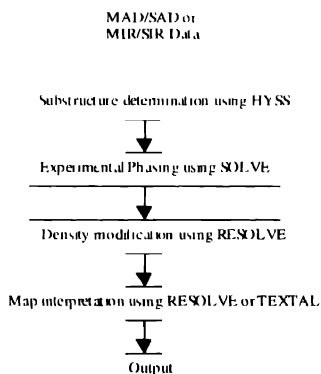
*Corresponding Author

has acquired increasing popularity and favorably competes with multiple wavelength anomalous diffraction (MAD). SAD may need more accurate intensity estimations than MAD, especially when weak anomalous scatterers such as phosphorous or sulfur are used [1,2] but does not require precise wavelength tuning and can be performed in home laboratories with Cu or Cr X-ray sources. The introduction of the chromium rotating anode and designs of new synchrotron beamlines especially for the use of long-wavelength radiation are deliberately intended for the application of SAD phasing [3].

PHENIX (Python-based Hierarchical Environment for Integrated Xtallography) is a software package developed for the automatic crystal structure determination. This provides the necessary algorithms to proceed from reduced intensity data to a refined molecular model [4,5].

2. Description of PHENIX program

PHENIX uses Python as the scripting language and C++ as the compiled language. High-level algorithms such as complex refinement protocols or phasing procedures can be developed in a scripting language and the computation of structure factors or discrete Fourier transforms must be implemented in a compiled language for performance reasons.



The algorithms implemented in PHENIX for automatic structure determination are shown in the flowchart. For a given data, Hybrid Substructure Search (HYSS) makes use of Patterson and direct methods to locate anomalous scatterers or heavy atoms for experimental phasing. Once the substructure has been determined, SOLVE program enables rapid configuration of jobs for experimental phasing through interfaces. Initial phases can also be obtained using molecular replacement, incorporating maximum likelihood targets. This can increase the success rate of this method using search models of lower structural similarity.

Maximum likelihood density modification algorithms implemented in RESOLVE produce minimally biased electron density maps using either the phases obtained from experimental phasing or by molecular replacement. These maps are then automatically interpreted using template matching implemented in RESOLVE and pattern recognition methods implemented in TEXTAL.

The strategy concept in PHENIX provides a way to connect complex networks of tasks to perform a higher-level task. For example, in SAD case, strategy concept breaks down steps required to go from initial data to a first electron density map in various tasks. The PHENIX graphical user interface permits strategies to be visualized and manipulated. The Data Storage (PDS) in PHENIX is a data management system that oversees the information generated for each step in structure determination and contains a complete history of each step in solution along with all of the generated structural information. The PHENIX GUI also provides PYMOL, a molecular graphics system written in C and Python, allowing easy viewing of structure and maps via close integration with PHENIX.

3. Materials and methods

As detailed papers on the success of SAD phasing using Cu K α radiation or synchrotron radiation with wavelength range 1.1–1.7 Å or Cr K α radiation have already appeared in the literature [6–11], this paper mainly focuses on the SAD application of PHENIX to an enzyme glucose isomerase (approximately 44 kDa molecular weight) using lab source (both CuK α and CrK α) anomalous scattering data corresponding to 1.7 Å.

Table 1 Crystallographic data

For CuK α data	
<i>a</i> (Å)	93.021
<i>b</i> (Å)	97.948
<i>c</i> (Å)	102.641
Space group	1222
Resolution range (Å)	10.1–7.1 (7.56–1.7)
Completeness (%)	99.91 (99.94)
<i>I</i> / σ (<i>I</i>)	28.38 (6.98)
Anomalous signal-to-noise ratio	0.99
For CrK α data	
<i>a</i> (Å)	92.872
<i>b</i> (Å)	97.905
<i>c</i> (Å)	102.710
Space group	1222
Resolution range (Å)	10.2–3.2 (3.77–2.3)
Completeness (%)	98.47 (99.06)
<i>I</i> / σ (<i>I</i>)	58.3 (42.22)
Anomalous signal-to-noise ratio	1.39

ditions. This enzyme contains 388 amino acids and two disulfide sites, one occupied by Mn^{2+} ion and the other by Mg^{2+} ion. Two data sets were collected at High Intensity X-ray diffractometer at Nagoya University using Rigaku R-axis IV and VII goniometer system. Table 1 shows the crystallographic details for these two data sets.

In cases where PHENIX did not build a complete model but a partial one, this model could be fed into the automatic model building or side chain tracing by the program ARP/wARP, followed by the refinement program REFMAC [13]. All the calculations mentioned here were carried out using Pentium 1

Glucose isomerase, $\text{CuK}\alpha$ data

(i) A manganese atom located from SAD data

PHENIX was run to find substructure and then for phasing and model building with the inputs of scalepack format for this data, a model letter sequence file and substructure present in this data as manganese (Mn). The imaginary component of the anomalous scattering (f'') of manganese at this wavelength is 1.6 electron units. Initially, ILYSS algorithm in PHENIX located manganese position (along with other heavy atoms) and model building was done using SOLVE algorithm with the manganese position. After density modification and map interpretation with the SOLVE algorithm in PHENIX, the program finally built residues in 21 chains out of a total of 388 residues. Out of these 21 residues, 31 residues were built with side chains. At this stage, PHENIX gave the best score of 25.8, Figure of Merit (MOS) of 0.17, R and R_f as 48% and 50%, respectively. The final model map correlation coefficient for this output was

0.4. A map was calculated using the SOLVE output phases and 492 peaks were above 3σ cut-off.

This model was then fed to ARP/wARP [12] for automatic model building using the option 'model building using the existing model'. After 50 cycles of auto-building, it was able to build 383 residues with side chains (out of a total of 388 residues) and has located 351 water atoms. At this stage, the R_w and R_f values were 17.1% and 20.6%, respectively. Without the water atoms R_w and R_f values were 23.3% and 24.7%, respectively. The map indicated the difference densities of the missing regions and the remaining residues were modeled into this. After the manual model building, the water atoms were checked and included if necessary and 25 cycles of maximum-likelihood refinement were performed using REFMAC [13]. The final R_w and R_f values were 17.3 and 19.4%, respectively. Table 2a details the PHENIX and ARP/WARP results.

(ii) Using 11 heavy atoms

Using the same data, PHENIX was run with the substructure containing eleven heavy atoms (9 sulphurs in Cys and Met residues + 1Mn + 1Mg). The imaginary component of the anomalous scattering (f'') of sulphur at this wavelength is 0.57 electron units. After the substructure finding, phasing and model building, PHENIX built 341 residues in one chain out of 388 residues. Out of these 341 residues, 335 residues were with side chains. The overall model map correlation coefficient for this output was 0.71. A map was calculated using the SOLVE output phases and 557 peaks were above 3σ cut-off. This model was then fed into ARP/wARP which built 382 residues with side chains and located 354 water atoms. Manual model building was carried out for the missing residues and the water atoms

Table 2a. Details of PHENIX and ARP/wARP glucose isomerase $\text{CuK}\alpha$ data, one Mn position, R_w and R_f are in %.

PROGRAM	Resolution limit	10-1.7 Å			
PHENIX	One peak	Score = 25.8 FOM = 0.17		SOLVE MCC = 0.2990 492 peaks	
	RESOLVE built	191 (31 with $R_w = 48\%$ side chains)		$R_f = 50\%$	Overall model MCC = 0.4
		Initial	$R_w = 47.6\%$ $R_f = 46.7\%$		
	No. of auto building cycles	10			
ARP/wARP	No. of REFMAC cycles in each auto building cycle	5			
		Final	$R_w = 17.1\%$ $R_f = 20.6\%$		
	Connectivity Index	0.99			
	No. Chains	2			
	No. Res. Built	383 (with side chains)			
	Water atoms	351			
	R_w and R_f without water atoms	$R_w = 23.3\%$ $R_f = 24.7\%$			
Final model with solvent atoms		$R_w = 17.3\%$ $R_f = 19.4\%$			
		r.m.s. deviation of backbone atoms (LOAD) 0.150 Å			
		(1MNZ) 0.269 Å			

were checked and included if necessary. Then 25 cycles of maximum-likelihood refinement were performed using REFMAC. The final R_w and R_f values were 17.7 and 20.1%, respectively. All these details are presented in Table 2b.

388 residues. The overall model map correlation coefficient; this output was 0.73. A map was calculated using the N_{obs} output phases and 285 peaks were above 3σ cut-off.

Table 2b. Details of PHENIX and ARP/wARP glucose isomerase CuK α data, 11 heavy atom positions, R_w and R_f are in %.

PROGRAM	Resolution limit	10.17 Å
PHENIX	Eleven peaks	Score = 63.0 FOM = 0.29 SOLVE MCC = 0.3626, 557 peaks
	RESOLVE built	341 (335 with $R_w = 31$ side chains) $R_f = 34$ Overall model MCC = 0.71
ARP/wARP		$R_w = 30.3$ $R_f = 30.4$
	No. of auto building cycles	10
	No. of Refmac cycles in each auto building cycle	5
	Connectivity Index	$R_w = 17.5$ $R_f = 20.9$
	No. Chains	0.98
	No. Res. Built	3
	Water atoms	382 (with side chains)
	R_w and R_f without water atoms	154
		$R_w = 23.5$ $R_f = 24.7$
Final model with solvent atoms		$R_w = 17.7$ $R_f = 20.1$
	r.m.s. deviation of backbone atoms(1 σ AD) 0.201 Å (1MNZ) 0.302 Å	

3.2 Glucose isomerase CrK α data

PHENIX was run with the inputs of scalepack format for this data set, single letter sequence file and substructure present in this data as sulphur (S). The imaginary component of the anomalous scattering (f'') of sulphur for this wavelength is 1.14 electron units. The phasing and model building procedure built 383 residues in two chains (256 residues with side chains) out of

To build the side chains of the remaining residues, this was led to ARP/wARP and finally it built 381 residues with side chains out of 388 residues and located 300 water atoms. Map model building was carried out for the missing residues and water atoms were checked and included if necessary. The cycles of maximum-likelihood refinement were performed using REFMAC. The final R_w and R_f values were 16.2 and 21.3%, respectively. All these details are presented in Table 3.

Table 3. Details of PHENIX and ARP/wARP glucose isomerase CrK α data, 11 heavy atom positions, R_w and R_f are in %.

	Resolution limit	10.23 Å
	Eleven peaks	Score = 71.8 FOM = 0.37 SOLVE MCC = 0.3452, 285 peaks
	RESOLVE built	383 (256 with $R_w = 32$ side chains) $R_f = 37$ Overall model MCC = 0.73
		$R_w = 31.9$ $R_f = 29.9$
	No. of auto building cycles	10
	No. of Refmac cycles in each auto building cycle	5
	Connectivity Index	$R_w = 16.7$ $R_f = 21.9$
	No. Chains	0.98
	No. Res. Built	3
	Water atoms	381 (with side chains)
	R_w and R_f without water atoms	300
		$R_w = 22.5$ $R_f = 25.7$
Final model with solvent atoms		$R_w = 16.2$ $R_f = 21.3$
	r.m.s. deviation of backbone atoms(1 σ AD) 0.204 Å (1MNZ) 0.373 Å	

Results and discussion

Glucose isomerase

single wavelength, manganese position

Figures 1a and 1b show the cartoon diagrams of the PHENIX model and the final model. The backbone of the final model is superimposed with glucose isomerase structures deposited in the protein data bank (PDB) (PDB id: 1OAD, space group $P2_1$; PDB id: 1MNZ; space group $I222$ – atomic resolution 1.1 Å). The root-mean-square deviation is 0.15 Å with PDB (1OAD) and 0.269 Å with PDB (1MNZ). Figure 1c shows a helix region of the final model superposed with SOLVE map and final $2F_o - 1F_c$ map. The map correlation coefficient between SOLVE map and the final map is 0.299. The average thermal factor of the current model is 14.9 Å² and the estimated overall coordinates error is 0.091 Å.

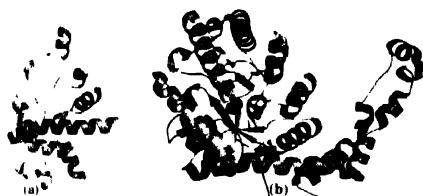


Figure 1a. PHENIX model (191 a.a. (31 a.a. with side chains) from one heavy atom for CuK α data.

Figure 1b. Final model (CuK α data) using one Mn Auto Built (383 a.a. with side chains) using ARP/wARP from PHENIX output.

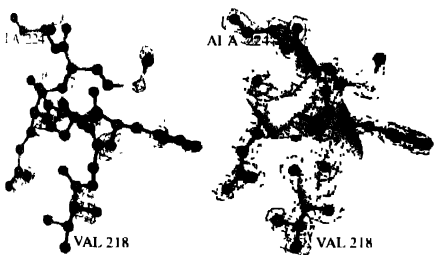


Figure 1c. Final model of helix region superposed with SOLVE map and final $2F_o - 1F_c$ map (0.9 σ).

heavy atoms

Figures 2a and 2b show the cartoon diagrams of the PHENIX model and the final model. The r.m.s. deviation of this model with the PDB id: 1OAD is 0.201 Å and that with PDB id: 1MNZ is 0.302 Å. Figure 2c shows a helix region of the final model superposed with SOLVE map and also final $2F_o - 1F_c$ map. The correlation coefficient between the SOLVE map and the final map is 0.3626. The average thermal factor of the current

model is 14.9 Å² and the estimated overall coordinates error is 0.096 Å.



Figure 2a. PHENIX model (341 a.a. (335 a.a. with side chains) from 11 heavy atoms for CuK α data.

Figure 2b. Final model (CuK α data) using 11 heavy atoms Auto Built (382 a.a. with side chains) using ARP/wARP from PHENIX output.

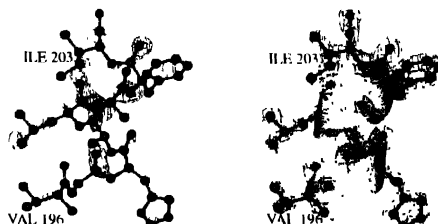


Figure 2c. Final model of helix region superposed with SOLVE map and final $2F_o - 1F_c$ map (0.9 σ).

4.2 Glucose isomerase, CrK α data

Figures 3a and 3b show the cartoon diagrams of the PHENIX model and the final model. The r.m.s. deviation of this model

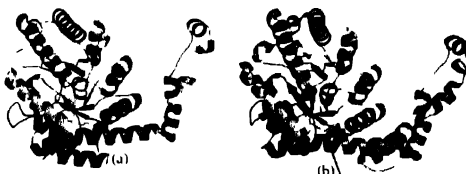


Figure 3a. PHENIX model (383 a.a. (256 a.a. with side chains) from 11 heavy atoms for CrK α data.

Figure 3b. Final model (CrK α data) using 11 heavy atoms Auto Built (381 a.a. with side chains) using ARP/wARP from PHENIX output.

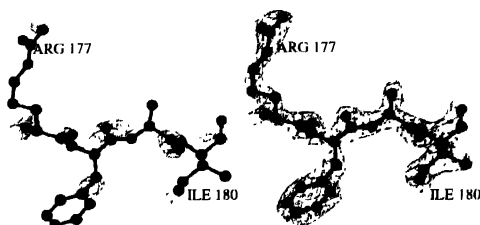


Figure 3c. Final model superposed with SOLVE map and final $2F_o - 1F_c$ map (0.7 σ).

with the PDB i.d. 1OAD is 0.204 Å and that with PDB i.d. 1MNZ is 0.373 Å. Figure 3c shows a section of the final model superposed with SOLVE map and also final $2F_o - |F_o|$ map. The map correlation coefficient between the SOLVE map and the final map is 0.3452. The average thermal factor of the current model is 13.8 Å² and the estimated overall coordinates error is 0.262 Å.

Figure 4a shows the superposition of the C α atoms of the current model with the final model of CuK α data obtained using one Mn, 11 heavy atoms, final model of CrK α data with PDB i.d. 1OAD. It shows that the overall fold of the models obtained using CuK α and CrK α data sets are similar to that of PDB i.d. 1OAD. Figure 4b shows the superposition of the current model with the final model using CuK α data set obtained with one Mn, 11 heavy atoms, final model of CrK α data and PDB 1MNZ.

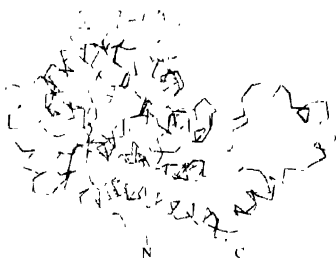


Figure 4a. Superposition of the C α atoms of the final model of CuK α data using 1Mn (blue), 11 heavy atoms (cyan), CrK α data using 11 heavy atoms (magenta) along with PDB 1OAD (red).



Figure 4b. Superposition of the C α atoms of the final model of CuK α data using 1Mn (blue), 11 heavy atoms (cyan), CrK α data using 11 heavy atoms (magenta) along with PDB 1MNZ (red).

5. Conclusion

As suggested by Wang [14], the longer wavelength from a chromium target is advantageous for sulphur SAD phasing because the f'' value of the S atom is larger (1.14 e⁻) than that obtained with copper radiation (0.57 e⁻). This can be confirmed from the results obtained for CrK α data, using eleven sulphur positions, PHENIX built only 341 residues for CuK α data whereas 383 residues were built for CrK α data. The above work

emphasizes the applicability of the SAD technique to solve macromolecular structure using lab source data using CrK α radiation when data extends to 2.3 Å resolution. SAD data of a macromolecule can be used to solve the structure with the existing sophisticated program PHENIX in an automated way using this approach. Automation not only speeds up the process of solution of crystal structures but also permits the selection of better and more accurate atomic models. Inspection of P statistics shows that the average quality of structures from the structural genomics centres does not differ from the quality of structures elucidated in a more traditional way, hence this method of solving a macromolecular structure using lab source data is beneficial to those who do not have access to synchrotron data collection.

Acknowledgment

SS thanks Council of Scientific and Industrial Research (CSIR) for providing Senior Research Fellowship. DV acknowledges Bioinformatics division of Department of Biotechnology (DBT) and University Grants Commission (UGC), Govt. of India for major projects supporting this work and thanks Venture Bioscience Laboratory authorities, Nagoya University, Nagoya, Japan for the visiting Professorship assignments in short terms and also acknowledges financial support to the Department under IIT-SAP and DST-FIST programmes.

- [1] Z. Dauter and D. A. Adams *Acta Cryst.* **D57** 990 (2001).
- [2] U. A. Ramagopal, M. Dauter and Z. Dauter *Acta Cryst.* **D59** (2003).
- [3] Z. Dauter *Acta Cryst.* **D62** 1 (2006).
- [4] P. D. Adams, R. W. Grosse-Kunstleve, L. Hung, T. R. Ioerger, M. McCoy, N. W. Moriarty, R. J. Read, J. C. Sacchettini, N. K. Sauter, T. C. Terwilliger *Acta Cryst.* **D58** 1948 (2002).
- [5] P. D. Adams, K. Gopal, R. W. Grosse-Kunstleve, L. Hung, T. R. Ioerger, M. McCoy, N. W. Moriarty, R. K. Pai, R. J. Read, T. D. Romo-Clavero, N. K. Sauter, L. C. Storoni and T. C. Terwilliger *Synchrotron Rad.* **11** 53 (2004).
- [6] C. S. Bond, M. P. Shaw, S. Alpheys and W. N. Hunter *Acta Cryst.* **D57** 755 (2001).
- [7] S. Selvanayagam, D. Velmurugan, T. Yamane and A. Suzuki *International Symposium on Recent Trends in Macromolecular Structure and Function Meeting, Poster P24* (2006).
- [8] E. J. Gordon, G. A. Leonard, S. McSweeney and P. F. Zagalski *Acta Cryst.* **D57** 1230 (2001).
- [9] K. Sekar, V. Rajakannan, D. Velmurugan, T. Yamane, M. Dauter and Z. Dauter *Acta Cryst.* **D60** 11 (2004).
- [10] Z. Dauter, M. Dauter, E. de La Fortelle, G. T. Wilson and G. Sheldrick *J. Mol. Biol.* **289** 83 (1999).
- [11] Y. Kitago, N. Watanabe and I. Tanaka *Acta Cryst.* **D61** 1011 (2005).
- [12] A. Perrakis, R. Morris and V. S. Lamzin *Nature Struct. Biol.* **6** 680 (1999).
- [13] G. N. Murshudov, A. Lebedev, A. A. Vagin, K. Wilson and J. Dodson *Acta Cryst.* **D55** 247 (1999).
- [14] B. C. Wang *Methods Enzymol.* **115** 90 (1987).